Phonemic Tagging with the Unisyn lexicon

LaBB-CAT includes the *Unisyn layer manager*, which is designed for ingesting Unisyn accent-specific lexicons. Unisyn must be downloaded separately, and the included scripts executed to produce a lexicon for the desired variety. The resulting file can be added to LaBB-CAT, and then the layer manager can be configured to use it for tagging word tokens with their phonemic transcriptions.

Unisyn is a 'master lexicon' of English, which contains:

- orthography
- part-of-speech
- pronunciation, in an 'accent neutral' form
- 'enriched orthography' showing morphological information
- frequency, as derived from various sources, including the British National Corpus, Time articles, Gutenberg, etc.

The pronunciations in the lexicon can be converted into an accent-specific form using perl scripts that are included with the lexicon.

Getting Unisyn

Unisyn is available under a non-commercial license, and must be acquired seperately from this layer manager. To acquire Unisyn, you must first register on the the Unisyn website and accept the terms of their license. The Unisyn website is here: http://www.cstr.ed.ac.uk/projects/unisyn/

(This layer manager has been tested with version 1.3 of Unisyn)

Install the Layer Manager

First, the Unisyn layer manager module must be installed:

- 1. Select the *layer managers* menu option.
- 2. Follow the *List of layer managers that are not yet installed* link near the bottom.
- 3. Find "Unisyn" in the list, and press its *Install* button, then *Install* again.

Once the layer manager is installed, you'll see a page of information explaining more details about the layer manager. Once you've read this information, you can close the browser tab to return to LaBB-CAT.

Using Unisyn with this layer manager

Once you've got Unisyn, you can use it to produce accent-specific lexicons, and provide these lexicons to the layer manager, which then uses them to annotate words in LaBB-CAT.

For example, if you want to annotate your transcripts with 'General American English' pronunciations:

- 1. Generate the General American English (*gam*) lexicon by running the following Unisyn commands:
 - I. get-exceptions.pl -a gam -f unilex > gam.1
 - 2. post-lex-rules.pl -a gam -f gam.1 > gam.2
 - 3. map-unique.pl -a gam -f gam.2 > gam.unisyn. This gives you the file gam.unisyn, which is the lexicon file you need for the next step.
- 2. Upload the accent-specific lexicon into LaBB-CAT:
 - 1. Select *layer managers* on the menu.
 - 2. Find the *Unisyn layer manager* in the list, and click the *Extensions* button (the second-to-last button on the right).
 - 3. Press *Choose File* and select the *gam.unisyn* file you generated above.
 - 4. Press *Upload Lexicon*. You will see a progress bar while the file is uploaded and the data is processed.
- 3. Create the layer for your pronunciation annotations: To create a new layer with CMUdict annotations:
 - 1. Select the *word layers* option on the menu this will display a list of all the word layers you already have in the database.
 - 2. At the top of the list, there's a blank form for creating a new layer fill this form in:
 - Layer ID: enter a name e.g. phonemes
 - **Type**: select *Phonological*
 - Manager: select Unisyn
 - Alignment: select *None* (as these are simply tags on the orthographic words)
 - Generate: select Always
 - **Description**: something like All possible pronnunciations according to Unisyn's General American accentspecific lexicon
 - 3. Press the *New* button to create the layer. You will see the layer configuration page. Check the online help for explanations of all options, but at least:
 - 4. Ensure the Source Layer is *orthography*
 - 5. Select the desired Lexicon from the list (these relate to the file or files you generated and uploaded above)

6. Tick the Strip syllabification/stress if you will use this layer for forced alignment

| Source Layer: | orthography 🗸 |
|-------------------------------|-----------------|
| Language: | en |
| Lexicon: | gam.unisyn 🗸 |
| Field: | Phonemes (DISC) |
| Strip syllabification/stress: | |
| First variant only: | |
| Recover Syllables: | |
| | Save |

- with HTK. 7. Press *Save*
- 8. Press *Regenerate*. You will see a progress bar while the layer manager annotates all the transcripts that have already been uploaded.

LaBB-CAT will then generate annotations for all the transcripts you already have in your database. If you have a lot of data, this may take a while.

From now on, when you upload a new transcript, the Unisyn annotations will automatically be generated for it.

Mapping Unisyn pronunciations to the DISC phoneme set

LaBB-CAT's processing of phonological layers assumes that the annotations use the DISC phoneme set designed for the CELEX phonemic transcriptions. This set is used because each phoneme is expressed by precisely one ASCII character, including phonemes usually expressed using a digraph - e.g. affricates like /t / (which is /J/ in DISC) and diphthongs like /a / (which is /2/ in DISC)

Unisyn transcriptions use a set of phones that is greater that the set of phones available in DISC, and the transcriptions are designed to be broadly phonetic, not phonemic.

This means that using the DISC representation of the transcripts is imperfect, as there is a certain amount of loss of information when mapping Unisyn phones to DISC phonemes. The default mapping that is used is shown below.

| Unisyn | | DISC | IPA | Lexical set e.g. |
|--------|---------------|------|-----|---------------------------------------------------------|
| ah | \rightarrow | # | e: | BATH |
| aa | \rightarrow | Q | D | $PALM \rightarrow LOT$ |
| ar | \rightarrow | Q | D | start \rightarrow PALM \rightarrow LOT |
| oa | \rightarrow | { | æ | $BANANA \rightarrow TRAP$ |
| ao | \rightarrow | # | a: | $MAZDA \rightarrow BATH$ |
| e | \rightarrow | Е | ε | DRESS |
| er | \rightarrow | Е | ε | r-coloured DRESS in scots en |
| а | \rightarrow | { | æ | TRAP |
| eh | \rightarrow | { | æ | ann use TRAP |
| ou | \rightarrow | 5 | ອບ | GOAT - but a monophthong for in some varieties |
| oul | \rightarrow | 5 | ອບ | goal - post vocalic GOAT |
| ouw | \rightarrow | 5 | əυ | KNOW \rightarrow GOAT (except for Abergave) |
| 0 | \rightarrow | Q | D | LOT |
| oou | \rightarrow | Q | D | $adios \rightarrow LOT$ |
| au | \rightarrow | Q | D | CLOTH \rightarrow LOT (but a diphthong in some en-US) |
| 00 | \rightarrow | \$ | э: | THOUGHT (but a diphthong in some varieties) |
| or | \rightarrow | \$ | э: | r-coloured THOUGHT |
| ii | \rightarrow | i | i: | FLEECE |
| iy | \rightarrow | i | i: | HAPPY - I for some varieties |
| ie | \rightarrow | i | i: | HARRIET - Leeds only |
| ii; | \rightarrow | i | i: | $AGREED \rightarrow FLEECE$ |
| ir | \rightarrow | i | i: | NEARING - r-coloured NEAR \rightarrow FLEECE |
| ir; | \rightarrow | i | i: | near - scots-long NEAR \rightarrow FLEECE |
| i | \rightarrow | Ι | Ι | KIT |
| @ | \rightarrow | @ | ə | schwa |
| @r | \rightarrow | @ | ə | r-coloured schwa |
| uh | \rightarrow | V | Λ | STRUT |
| u | \rightarrow | U | υ | FOOT |
| uu | \rightarrow | u | u: | GOOSE |
| iu | \rightarrow | u | u: | $BLEW \rightarrow GOOSE$ |
| uu; | \rightarrow | u | u: | brewed \rightarrow GOOSE |
| uw | \rightarrow | u | u: | louise \rightarrow GOOSE |
| uul | \rightarrow | u | u: | goul - post-vocalic GOOSE |
| ei | \rightarrow | Ι | еі | FACE |
| ee | \rightarrow | Ι | еі | WASTE \rightarrow FACE (except for abercrave) |
| ai | \rightarrow | 2 | аі | PRICE |
| ae | \rightarrow | 2 | аі | TIED \rightarrow PRICE (except Edi and Aberdeen) |
| ae | \rightarrow | 2 | аі | TIED \rightarrow PRICE (except Edi and Aberdeen) |
| aer | \rightarrow | 2 | аі | FIRE - r-coloured PRICE |
| aai | \rightarrow | 2 | аі | TIME \Rightarrow PRICE (except S. Carolina) |
| oir | \rightarrow | 2 | аі | COIR - r-coloured PRICE |

| Unisyn | | DISC | IPA | Lexical set e.g. |
|--------|---------------|------|-----|----------------------------------------------------------|
| @@r | \rightarrow | 3 | 3: | NURSE |
| oi | \rightarrow | 4 | ЭI | CHOICE |
| ow | \rightarrow | 6 | au | MOUTH |
| owr | \rightarrow | 6 | au | HOUR - r-coloured MOUTH |
| oow | \rightarrow | 6 | au | HOUR \rightarrow MOUTH (exception S. Carolina) |
| i@ | \rightarrow | 7 | IÐ | NEAR |
| iir | \rightarrow | 7 | IÐ | beard \rightarrow NEAR (except en-AU) |
| eir | \rightarrow | 8 | 63 | SQUARING (actually a monophthong in many varieties) |
| ur | \rightarrow | 9 | υə | JURY |
| ur; | \rightarrow | 9 | ΰə | CURE - scots-long JURY |
| iur | \rightarrow | 9 | υə | curious - JURY exception in Cardiff & Abercrave |
| р | \rightarrow | р | р | |
| t | \rightarrow | t | t | |
| ? | \rightarrow | ? | ? | (glottal stop) |
| t^ | \rightarrow | L | ſ | butter/merry flap |
| k | \rightarrow | k | k | |
| X | \rightarrow | Х | Х | loch |
| b | \rightarrow | b | b | |
| d | \rightarrow | d | d | |
| g | \rightarrow | g | g | |
| ch | \rightarrow | J | ţ | |
| jh | \rightarrow | _ | dз | |
| S | \rightarrow | S | S | |
| Z | \rightarrow | Z | Z | |
| sh | \rightarrow | S | ſ | |
| zh | \rightarrow | Z | 3 | |
| t | \rightarrow | f | f | |
| V | \rightarrow | V | V | |
| th | \rightarrow | T | θ | |
| dh | \rightarrow | D | ð | |
| h | \rightarrow | m | m | |
| m | \rightarrow | m | m | |
| m! | \rightarrow | F | m | chasm |
| n | \rightarrow | n | n | |
| n! | \rightarrow | H | ņ | mission |
| ng | \rightarrow | N | ŋ | |
| 1 | \rightarrow | 1 | 1 | |
| 11 | \rightarrow | 1 | 1 | llandudno (for Cardiff and Abercrave, this is different) |
| lw | \rightarrow | 1 | 1 | teel - dark l |
| r | \rightarrow | r | r | |
| У | \rightarrow | j | j | |

| Unisyn | | DISC | IPA | Lexical set e.g. |
|--------|---------------|------|-----|------------------|
| w | \rightarrow | W | W | which |
| hw | \rightarrow | w | W | which |

Changing the DISC symbol mapping

If you would like to adjust this mapping, so that Unisyn symbols correspond to different DISC symbols, you can do so by visitin the Unisyn Layer Manager's *Extensions* page, and following the dictionary's 'Phoneme Map' link, as shown in Figure 1

Unisyn Tagger

The Unisyn Tagger tags word tokens with data from <u>Unisyn</u>, a lexicon produced by the Centre for Speech Technology Research at the University of Edinburgh.

On this page, you can add or delete lexicons.

Lexicons

gam.unisyn Phoneme Map X

gnz.unisyn Phoneme Map X

• Upload New Lexicon File

Figure 1: Lexicon listing, where mappings can be speccified, and lexicons uploaded, or deleted

This page allows you to change how Unisyn symbols correspond to DISC symbols, as seen in Figure 2. You can change symbols/text in the text boxes, and then save your changes with the save button that appears at the bottom of the list.

Other possible phoneme encodings

If having the original transcriptions precisely as defined in the Unisyn lexicon is very important, you can instead create a layer that uses the original transcription as contained in the file you uploaded. This has the advantage that the transcriptions are not filtered through the above mapping, and the disadvantage that LaBB-CAT won't be able to display the transcriptions using IPA symbols, nor help you when creating search patterns for the layer.

If you decide to do this, Unisyn offers you two possible representations:

• Unisyn transcriptions - e.g. { p r @ . n ~ uh n s \$}.> ii . * ei . sh n! > - these are already present in the file that you generated if you followed the instructions above (i.e. gam.unisyn)

Unisyn Tagger gam.unisyn

| Original | → DISC | Note | + |
|----------|--------|----------------------------------|---|
| * | • | primary stress | X |
| - | , | tertiary stress | X |
| | - | syllable boundary | X |
| @ | @ | schwa | X |
| @@r | 3 | NURSE | X |
| @r | @ | r-coloured schwa | X |
| а | { | TRAP | X |
| аа | Q | PALM -> LOT (US) (but BATH (RI | X |
| aai | 2 | TIME -> PRICE (except S. Carolin | X |
| ae | 2 | TIED -> PRICE (except Edi and Al | X |
| aer | 2 | FIRE - r-coloured PRICE | X |
| ah | # | BATH | X |
| ai | 2 | PRICE | X |

()Phoneme Map

Figure 2: Phoneme Map page allowing correspondences to be edited

• SAM-PA transcriptions - e.g. pr\@%nVns\$i"e\$Sn=\$@5 - these can be obtained by running an extra Unisyn command, and uploading the resulting gam.sampa file: output-sam.pl -a gam -f gam.unisyn > gam.sampa

(Unisyn has a third script called output-ipa.pl which produces transcriptions for displaying in HTML - e.g. p \Rightarrow n ns.i \in . n - which are not suitable for search, analysis, or forcedalignment)

In order to prevent the DISC mapping from applying on your layer:

- When creating the layer, set the layer type to *Text* rather than *Phonological*.
- When configuring the layer, set the field to *Phonemes (original file)* rather than *Phonemes (DISC)*.